

AD-A104 480 JOHNS HOPKINS UNIV BALTIMORE MD DEPT OF MATHEMATICAL--ETC F 6 12/1
GENERALIZED L-, M- AND R-STATISTICS. (U) AUG 81 R J SERFLING N00014-79-C-0801
UNCLASSIFIED TR-349

UNCLASSIFIED TR-349

N00014-79-C-0801

1 OF 1
AD A
- (AASHI)

END
DATE
FILED
10-81
DTIC

DEPARTMENT OF MATHEMATICAL SCIENCES
The Johns Hopkins University
Baltimore, Maryland 21218

LEVEL

GENERALIZED L-, M- AND R-STATISTICS,

by

Robert J. Serfling
The Johns Hopkins University

Technical Report No. 349
ONR Technical Report No. 81-8

August, 1981

DTIC
SELECTED
SEP 23 1981

Research supported by the Army, Navy and Air Force under Office
of Naval Research Contract No. N00014-79-C-0861. Reproduction
in whole or in part is permitted for any purpose of the
United States Government.

NT A

AD A104480

FILE COPY

ABSTRACT

GENERALIZED L-, M- AND R-STATISTICS

A class of statistics generalizing U-statistics and L-statistics, and containing other varieties of statistic as well, such as trimmed U-statistics, is studied. Using the differentiable statistical function approach, differential approximations are obtained and the influence curves of these generalized L-statistics are derived. These results are employed to establish asymptotic normality for such statistics. Parallel generalizations of M- and R-statistics are noted. Strong convergence, Berry-Esséen rates, and computational aspects are discussed.

1. Introduction. We consider a new class of statistics, which usefully generalizes the classes of U-statistics and L-statistics and contains other varieties of statistic as well. Let X_1, \dots, X_n be independent random variables having common probability distribution F . (More generally, the X_i 's could be random elements of an arbitrary space.) Let a "kernel" $h(x_1, \dots, x_m)$ be given, which for convenience and without loss of generality is assumed to be symmetric in its arguments. Denote by

$$(1.1) \quad W_{n,1} \leq \dots \leq W_n,$$

the ordered values of $h(X_{i_1}, \dots, X_{i_m})$ taken over the $\binom{n}{m}$ subsets of m distinct elements i_1, \dots, i_m from $\{1, \dots, n\}$. Consider the statistics given by

$$(1.2) \quad \sum_{i=1}^m c_{n,i} W_{n,i},$$

where $c_{n,i}$, $1 \leq i \leq \binom{n}{m}$, are arbitrary constants. The form (1.2) is quite general. It includes the U-statistic corresponding to the kernel h , which is given by (1.2) with $c_{n,i} = 1/\binom{n}{m}$, all i . And it includes the class of L-statistics (linear functions of order statistics), given by (1.2) for the particular kernel $h(x) = x$. Moreover, it includes statistics such as

$$\text{median } \{i(X_{i_1} + X_{i_2}), 1 \leq i\},$$

which is a standard version of the well-known Hodges-Lehmann location estimator.

Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/	
Availability Codes	
Avail and/or	
Dist	Special

R

AS 1970 subject classifications: Primary 62E20, Secondary 60F05

Key words and phrases: order statistics; L-statistics; M-statistics; R-statistics; Hodges-Lehmann estimator; trimmed U-statistics; asymptotic normality.

but which is neither a U-statistic nor an L-statistic. Thus, for example, the sample mean, the sample median (a particular L-statistic), the sample variance (a particular U-statistic), the Hodges-Lehmann location estimator, and the 5% trimmed mean (an L-statistic) - a group of statistics which traditionally have been viewed and analyzed as quite different types - may in fact be viewed from a single standpoint. In this way the form (1.2) provides a unifying concept relative to various familiar statistics. But (1.2) also embraces important new varieties of statistic. For example, "trimmed U-statistics" and "Winsorized U-statistics" fall in this class. In particular, a "trimmed variance" is defined by trimming the U-statistic corresponding to the kernel $h(x_1, x_2) = 4(x_1 - x_2)^2$. This provides a competitor to a somewhat similar nonparametric dispersion measure of Bickel and Lehmann (1976). Their measure is a trimmed variance which is simpler in form than a trimmed U-statistic but which is constructed assuming that the population median is known and incorporating its value into the measure.

Computationally, statistics requiring ordering such as the sample median and the Hodges-Lehmann have been deemed less satisfactory than statistics computed by averaging (such as the sample mean, the U-statistics) or by solving equations (e.g., M-estimates). It has appeared, and been asserted, that computation of the sample median required $O(n \log n)$ operations and that computation of the Hodges-Lehmann required $O(n^2)$ operations. However, with the advent of computer science and the development of ingenious algorithms for machine computation, this misunderstanding has been corrected. Indeed, the sample median requires only $O(n)$ operations (see Blum et al. (1973) and Floyd and Rivest (1975)) and the Hodges-Lehmann only $O(n \log n)$ operations (see Sharpen (1976)). Therefore, statistics of form (1.2) are not necessarily more formidable for machine computation than simpler types of statistic.

Despite the complexity and generality of the form (1.2), the usual asymptotic normality and convergence properties hold and can be expressed in explicit form. It turns out that for theoretical study of the class (1.2) it is appropriate to view the class as a generalization of L-statistics - hence the terminology "generalized L-statistics." This will become evident from the developments of Sections 2 and 3. In Section 2 statistics of form (1.2) will be represented as "statistical functions," i.e., as functionals of an empirical distribution function, in the spirit of von Mises (1947). The corresponding differential approximations will be derived, leading also to the influence curves of Hampel (1974). Results on asymptotic normality will be obtained in Section 3, using the results in Section 2 in conjunction with von Mises' approach of differentiable statistical functions, through a development parallel to the treatment of L-statistics in Serfling (1980), Chapter 8.

By analogy with the development of Section 2, whereby "generalized L-statistics" are formulated as statistical functions, one can formulate generalized M-statistics and R-statistics. These and other complements are discussed in Section 4.

An interesting feature of the treatment of generalized L-, M- and R-statistics is that the role played by the usual sample distribution function in the treatment of simple L-, M- and R-statistics is given over to a more complicated type of empirical distribution function, one having the structure of a general U-statistic. Accordingly, interesting generalizations of the well-developed theory of the usual empirical process become needed.

2. Generalized L-statistics: formulation, differential approximations,

and influence curves. The representation of a statistic as a functional, evaluated at a sample distribution function which estimates the underlying serial distribution function, helps to identify what parameter the statistic in question is actually estimating. It also sets the stage for application of differentiation methodology and influence curve analysis. Let us examine generalized L-statistics relative to these aims. We proceed by analogy with the treatment of simple L-statistics.

As before, we consider a sample X_1, \dots, X_n of independent observations having distribution F . Denote by F_n the usual sample distribution function,

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x), \quad -\infty < x < \infty,$$

where $I(A) = 1$ or 0 according as the event A holds or not. The class of (simple) L-statistics may be represented in the form

$$(2.1) \quad \sum_{i=1}^n c_{n,i} F_n^{-1}(i/n),$$

of which a suitably wide subclass can be represented as $T(F_n)$ for a functional $T(\cdot)$ of the form

$$(2.2) \quad T(F) = \int_0^1 F^{-1}(t) J(t) dt + \sum_{j=1}^d a_j F^{-1}(p_j).$$

Such a functional weights the quantiles $F^{-1}(t)$, $0 < t < 1$, of F according to a specified function $J(\cdot)$ for smooth weighting and/or specified weights a_1, \dots, a_d for discrete weighting. A particular L-functional is thus determined by specifying $J(\cdot)$, d , p_1, \dots, p_d and a_1, \dots, a_d . The corresponding L-statistic is then simply $T(F_n)$. Note that $T(F_n)$ may be written in the form

$$(2.3) \quad T(F_n) = \sum_{i=1}^n \left(\frac{i}{n} - \frac{1}{n} \right) J(i/n) + \sum_{j=1}^d a_j F_n^{-1}(p_j),$$

which exhibits the statistic explicitly as a linear function of the order statistics $F_n^{-1}(1/n)$, $1 \leq n$.

We now designate an analogous subclass of the statistics of form (1.2). For a given kernel $h(x_1, \dots, x_m)$, let H_h denote the empirical distribution function of the evaluations $h(X_{i_1}, \dots, X_{i_m})$, i.e.,

$$H_h(y) = \frac{1}{(n)} \sum_{i=1}^n I(h(X_{i_1}, \dots, X_{i_m}) \leq y), \quad -\infty < y < \infty,$$

where the sum is taken over the $\binom{n}{m}$ combinations of m distinct elements $\{i_1, \dots, i_m\}$ from $\{1, \dots, n\}$. The statistics of form (1.2) may be represented in the form

$$\sum_{i=1}^{\binom{n}{m}} c_{n,i} H_h^{-1}(i/\binom{n}{m}),$$

and, by analogy with (2.1) and (2.3), a wide and useful subclass of (1.2) is thus given in terms of the functional (2.2), by

$$(2.4) \quad T(H_h) = \sum_{i=1}^{\binom{n}{m}} \left(\frac{i}{(n)} - \frac{1}{n} \right) J(i/n) / \binom{n}{m} J'(i/n) + \sum_{j=1}^d a_j H_h^{-1}(p_j).$$

The parameter estimated by this "generalized L-statistic" (GL-statistic) is given by $T(H_p)$, where H_p is the distribution function estimated by H_h , i.e., where

$$H_p(y) = P_F(h(X_1, \dots, X_m) \leq y), \quad -\infty < y < \infty,$$

the distribution function of the random variable $h(X_1, \dots, X_m)$.

This functional approach allows the estimation error $T(H_n) - T(H_p)$ to be approximated by a differential quantity, which in practice can be obtained as a certain Gâteaux differential. As in Serfling (1980), Chapter 6, let us in general define the k -th order Gâteaux differential of a functional T at a distribution F in the direction of a distribution G to be

$$(2.5) \quad d_k T(F; G - F) = \frac{d^k}{d\lambda^k} T(F + \lambda(G - F))|_{\lambda=0} .$$

provided that the given right-hand derivative exists. For the simple L-functional $T(\cdot)$ given by (2.2), we have the first-order Gâteaux differential

$$(2.6) \quad d_1 T(F; G - F) = - \int_{-\infty}^{\infty} (G(y) - F(y)) J(F(y)) dy + \sum_{j=1}^d a_j \frac{p_j - G^{-1}(p_j)}{f(G^{-1}(p_j))} ,$$

assuming that F has a positive density f in neighborhoods of p_1, \dots, p_d (see Huber (1977) or Serfling (1980) for details). Accordingly, the estimation error $T(H_n) - T(H_p)$ of a GL-statistic becomes approximated by

$$(2.7) \quad d_1 T(H_n; H_p - H_n) = - \int_{-\infty}^{\infty} (H_n(y) - H_p(y)) J(H_p(y)) dy$$

$$+ \sum_{j=1}^d a_j \frac{p_j - H_n^{-1}(p_j)}{h_p(H_n^{-1}(p_j))} .$$

where h_p denotes the density of H_p , assumed to exist and be positive at p_1, \dots, p_d . A basic difference between the treatment of simple and generalized L-statistics, even though the same functional $T(\cdot)$ is involved in both cases, is that the quantity in (2.7) is a U-statistic in the more general case, but simply an average of IID's in the simple case. This stems from the fact that

H_n , which assumes in the general treatment the role played by F in the simple case, is a U-statistic. That is, for each fixed y , $H_n(y)$ is the U-statistic corresponding to the kernel $I(h(x_1, \dots, x_n) \leq y)$. Consequently, $d_1 T(H_p; H_n - H_p)$ is seen to be the U-statistic corresponding to the kernel

$$(2.8) \quad A(x_1, \dots, x_n) = - \int_{-\infty}^{\infty} (I(h(x_1, \dots, x_n) \leq y) - h_p(y)) J(H_p(y)) dy$$

$$+ \sum_{j=1}^d a_j \frac{p_j - I(h(x_1, \dots, x_n) \leq p_j)}{h_p(H_p^{-1}(p_j))} \leq h_p^{-1}(p_j) .$$

The formulas (2.7) and (2.8) will be relevant in treating the convergence theory of $T(H_n)$ in Section 3.

Also, formula (2.8) may be interpreted as an analogue of the usual influence curve. In the special case of a simple L-statistic, the "influence curve" associated with the statistic $T(H_n) = T(F_n)$ is obtained by putting $G = \delta_x$ (the distribution placing mass 1 at x) in the formula (2.6), which then yields the function $A(x)$ given by (2.8) with $h(x) = x$. In this case $A(X_1)$ represents the approximate "influence" of the observation X_1 on the estimation error when $T(F)$ is estimated by $T(F_n)$. (This interpretation, due to Hampel (1968), has become a standard concept in robustness considerations; see also Hampel (1974), Huber (1977).) Proceeding now to the generalized L-statistic, we see that $A(X_1, \dots, X_m)$ may be interpreted as the approximate influence of the combination of observations X_1, \dots, X_m on the estimation error when $T(H_p)$ is estimated by $T(H_n)$.

First, consider the special case of the functional

$$\tilde{T}_0(F) = H_F^{-1}(p).$$

When the parameter of interest is represented by $T(H_F)$, for some functional T evaluated at a distribution H_F related to the distribution F of the observations.

It is natural to use the estimator $T(H_n)$ based on an estimator H_n of H_F . As we have seen, however, the fact that H_n is in general a U-statistic introduces complications not present in the case of simple L-statistics, $T(F_n)$. Therefore, it is of some interest to view the parameter $T(H_F)$ as also, equivalently, the evaluation of some functional \tilde{T} at the basic distribution F . That is, $\tilde{T}(\cdot)$ is defined by

$$(2.9) \quad \tilde{T}(F) = T(H_F).$$

From this standpoint, a natural estimator is $\tilde{T}(F_n)$, or equivalently $T(H_n)$, where by definition

$$(2.10) \quad \begin{aligned} H_n(y) &= \int \dots \int I(y \leq x_1, \dots, x_m) \leq y dF_n(x_1) \dots dF_n(x_m) \\ &= \frac{1}{n^m} \sum_{i_1=1}^n \dots \sum_{i_m=1}^n I(y \leq X_{i_1}, \dots, X_{i_m}) \leq y. \end{aligned}$$

Note that $H_n(y)$ and $H_n^{-1}(y)$ are somewhat different, although closely related estimators. Thus $T(H_n)$ and $\tilde{T}(F_n) = T(H_n)$ are two somewhat different estimators of the "single parameter" expressed in two ways by (2.9). Although H_n is less straightforward than H_F for estimation of H_F , the estimator $T(H_n)$ is less

straightforward than H_F for estimation of H_F , the estimator $T(H_n)$ lends itself more straightforwardly to a standard influence curve analysis. Therefore, we derive the Gâteaux derivative of the functional \tilde{T} , as follows.

where $0 < p < 1$, and assume that H_F has a density h_F in the neighborhood of $H_F^{-1}(p)$, with $h_F(H_F^{-1}(p)) > 0$. Put

$$F_\lambda = F + \lambda(G - F)$$

and write

$$(2.11) \quad H_{F_\lambda}(H_{F_\lambda}^{-1}(p)) = p.$$

We may differentiate implicitly in (2.11) and solve for

$$\frac{d}{d\lambda} H_{F_\lambda}^{-1}(p)|_{\lambda=0+} = \frac{d}{d\lambda} \tilde{T}_0(F_\lambda)|_{\lambda=0+} = d_1 \tilde{T}_0(F; G - F).$$

For this purpose, we write

$$H_{F_\lambda}(y) = \sum_{j=0}^m \binom{m}{j} \lambda^{m-j} q_{F, G, j}(y),$$

where

$$q_{F, G, j}(y) = \int \dots \int I(m(x_1, \dots, x_m) \leq y) \prod_{i=1}^m dF(x_i) - F(x_j).$$

Then (2.11) becomes

$$(2.12) \quad \sum_{j=0}^m \binom{m}{j} \lambda^{m-j} q_{F, G, j}(H_{F_\lambda}^{-1}(p)) = p.$$

Now differentiate with respect to λ , obtaining the equation

$$(2.13) \quad \begin{aligned} &\sum_{j=0}^m \binom{m}{j} (m-j) \lambda^{m-j-1} q_{F, G, j}(H_{F_\lambda}^{-1}(p)) + \\ &+ \sum_{j=0}^m \binom{m}{j} \lambda^{m-j} q_{F, G, j}(H_{F_\lambda}^{-1}(p)) \frac{d}{d\lambda} H_{F_\lambda}^{-1}(p) = 0. \end{aligned}$$

where

$$q_{F,G,J}(y) = \frac{d}{dy} Q_{F,G,J}(y).$$

Letting $\lambda \rightarrow 0+$ in (2.13) and solving the resulting equation, we obtain

$$(2.14) \quad d_1 \tilde{T}_0(F;G,F) = -\frac{Q_{F,G,m-1}(\tilde{T}_0(F))}{q_{F,G,m}(\tilde{T}_0(F))},$$

and thus

(For higher-order Gâteaux derivatives of \tilde{T}_0 , we simply differentiate repeatedly in (2.12), or (2.13).) Since

$$Q_{F,G,m}(y) = h_F(y) = H_F(y),$$

and thus

$$Q_{F,G,m-1}(y) = \int \dots \int h(x_1, \dots, x_m) \leq y \prod_{i=1}^{m-1} dF(x_i) dG(x_{i+1}) = H_F(y),$$

and since also

$$(2.14) \quad \text{becomes} \quad (2.15) \quad d_1 \tilde{T}_0(F;G,F) = \frac{P - \int \dots \int h(x_1, \dots, x_m) \leq H_F^{-1}(p) \prod_{i=1}^{m-1} dF(x_i) dG(x_{i+1})}{h_F(H_F^{-1}(p))}.$$

In particular, putting $G=\delta_x$ (the distribution placing mass 1 at x), we obtain the influence curve of the functional \tilde{T}_0 ,

$$(2.16) \quad IC(x; \tilde{T}_0, F) = \frac{P - \int \dots \int h(x_1, \dots, x_{m-1}, x) \leq H_F^{-1}(p) \prod_{i=1}^{m-1} dF(x_i)}{h_F(H_F^{-1}(p))}$$

The following examples briefly illustrate these results for some familiar cases.

EXAMPLE 2.1: the p-th quantile of F. Take $m=1$, $h(x)=x$. Then

$$H_F(y) = P_F(X \leq y) = F(y),$$

$$\tilde{T}_0(F) = F^{-1}(p),$$

$$J_1(x; F^{-1}(p)) dG(x) = G(F^{-1}(p)),$$

and thus

$$(2.17) \quad d_1 \tilde{T}_0(F; G, F) = \frac{P - G(F^{-1}(p))}{f(F^{-1}(p))},$$

where f denotes the (assumed) density of F . Further, putting $G=\delta_x$, we obtain

$$(2.18) \quad IC(x; F^{-1}(p), F) = \frac{P - I(F^{-1}(p); 2x)}{f(F^{-1}(p))}.$$

These are well-known formulas for the p-th quantile functional. \square

EXAMPLE 2.2: Hodges-Lehmann type functionals. Take $m=2$ and

$$h(x_1, x_2) = \frac{1}{2}(x_1 + x_2). \quad \text{Then}$$

$$H_F(y) = P_F(\{x_1 + x_2 \leq y\}),$$

$$\tilde{T}_0(F) = H_F^{-1}(p),$$

$$J_1(x_1, x_2; H_F^{-1}(p)) dF(x_1) dG(x_2) = f(F(2H_F^{-1}(p) - x)) dG(x).$$

and thus

$$d_1 \tilde{T}_0(F; G, F) = 2 \frac{P - f(F(2H_F^{-1}(p) - x)) dG(x)}{h_F(H_F^{-1}(p))}.$$

In particular, for $G=\delta_x$, we have

$$IC(x; H_F^{-1}(p), F) = 2 \frac{P - f(F(2H_F^{-1}(p) - x))}{h_F(H_F^{-1}(p))}.$$

For the case F_n and F continuous and symmetric about 0, this yields

$$IC(x; H_F^{-1}(x), F) = 2 \frac{F(x)^{-1}}{\int F(x) dx}, \quad x > 0.$$

the familiar influence curve of the usual Hodges-Lehmann estimator. \square

Let us now consider the more general functional \tilde{T} given by (2.9), i.e.,

$$(2.19) \quad \tilde{T}(F) = \int_0^1 h_F^{-1}(t) J(t) dt + \int_{j=1}^d a_j h_F^{-1}(p_j).$$

Assuming validity of the interchange of order of differentiation and integration, in

$$\frac{d}{dt} \int_0^1 h_F^{-1}(t) J(t) dt = \int_0^1 \frac{d}{dt} h_F^{-1}(t) J(t) dt,$$

it follows in straightforward fashion, from the results for $\tilde{T}_0(\cdot)$, that

$$(2.20) \quad d_T(T(F; G; F)) = -m \sum_{j=1}^m \left\{ \int \dots \int I(h(x_1, \dots, x_m) \leq y) \prod_{i=1}^{m-1} dF(x_i) dG(x_m) - h_F(y) J(H_F(y)) \right\} dy \\ + m \sum_{j=1}^d a_j \frac{p_j - \int I(h(x_1, \dots, x_m) \leq h_F^{-1}(p_j)) \prod_{i=1}^{m-1} dF(x_i) dG(x_m)}{h_F(h_F^{-1}(p_j))}.$$

Accordingly, the influence curve of the functional $\tilde{T}(\cdot)$ in (2.19) is

$$(2.21) \quad IC(x; T, F) = -m \sum_{j=1}^m \left\{ \int \dots \int I(h(x_1, \dots, x_{m-1}, x) \leq y) \prod_{i=1}^{m-1} dF(x_i) - h_F(y) J(H_F(y)) \right\} dy \\ + m \sum_{j=1}^d a_j \frac{p_j - \int I(h(x_1, \dots, x_{m-1}, x) \leq h_F^{-1}(p_j)) \prod_{i=1}^{m-1} dF(x_i) - h_F(y) J(H_F(y))}{h_F(h_F^{-1}(p_j))}.$$

Since $\tilde{T}(F_n) - \tilde{T}(F)$ is approximately (under appropriate conditions)

$$(2.22) \quad d_1 \tilde{T}(F; F_n; F) = \frac{1}{n} \sum_{i=1}^n IC(X_i; \tilde{T}, F).$$

the IC represents the approximate "influence" of the observation X_i on the estimation error, when $\tilde{T}(F)$ is estimated by $\tilde{T}(F_n)$.

3. Asymptotic normality of GL-statistics. Under appropriate conditions

the GL-statistics $T(H_N)$ and $\tilde{T}(F_N)$ are asymptotically normal in distribution:

$$(3.1) \quad n^{1/2} (T(H_N) - T(H_F)) \stackrel{d}{\rightarrow} N(0, \sigma^2(T, H_F)).$$

and

$$(3.2) \quad n^{1/2} (\tilde{T}(F_N) - \tilde{T}(F)) \stackrel{d}{\rightarrow} N(0, \sigma^2(\tilde{T}, F)),$$

where $\sigma^2(T, H_F) = \sigma^2(\tilde{T}, F) = \sigma^2$ is given by

$$(3.3) \quad \sigma^2 = Var(IC(X; \tilde{T}, F)).$$

(Here $T(\cdot)$, $\tilde{T}(\cdot)$ and $IC(\cdot)$ are as defined in (2.2), (2.9) and (2.21). The asymptotic normality of $T(H_N)$ is established by making respectively.) The asymptotic normality of $\tilde{T}(H_N)$ by $d_1 T(H_N) - T(H_F)$ by $d_1 T(H_N) - \tilde{T}(H_F)$ as given in (2.7), rigorous the approximation of $\tilde{T}(H_N) - T(H_F)$ by $d_1 T(H_N) - \tilde{T}(H_F)$,

in which case (3.1) follows immediately from U-statistic theory (e.g., Serfling (1980, Chapter 5) and the asymptotic variance $\sigma^2(T, H_F)$ is given by

$$Var(A_1(X)), \text{ where } A_1(x) = E(A(x, X_1, \dots, X_{m-1})) \text{ and } A(x_1, \dots, x_m) \text{ is the function}$$

in (2.8). However, it is readily seen that $A_1(x) = IC(x; \tilde{T}, F)$, so that (3.3)

is valid. Likewise, the asymptotic normality of $\tilde{T}(F_N)$ is established by approximating $\tilde{T}(F_N) - \tilde{T}(F)$ by $d_1 \tilde{T}(F; F_n; F)$ and utilizing (2.22), in which case

(3.2) follows directly from classical central limit theory and the appropriate asymptotic variance is given immediately by (3.3). Specifically, these assertions are formalized in the following results (we shall deal explicitly only with $T(H_N)$ and discuss $\tilde{T}(F_N)$ in Remark 3.2(ii) at the end of this section).

THEOREM 3.1. Let H_F have positive derivatives at its p_j -quantiles. L.s.m. Let $J(t)$ vanish for t outside $[a, b]$, where $0 < a < b < 1$, and suppose that on (a, b) , J is bounded and continuous a.e. Lebesgue and a.e. H_F^{-1} . Assume that $0 < \sigma^2(T, H_F) < \infty$. Then (3.1) holds.

The foregoing result applies to examples such as trimmed and Winsorized U-statistics. The following result applies to untrimmed J functions.

THEOREM 3.2. Let H_F satisfy $\int (H_F(y)(1-H_F(y)))^{1/2} dy < \infty$ and have positive derivatives at its p_j -quantiles, L.s.m. Let J be continuous on $[0, 1]$. Assume that $0 < \sigma^2(T, H_F) < \infty$. Then (3.1) holds.

To prove these results, the functional $T(H_F)$ is treated in two parts,

$$T(H_F) = T_1(H_F) + T_2(H_F), \text{ where}$$

$$T_1(H_F) = \int_0^1 J(t) H_F^{-1}(t) dt$$

and

$$T_2(H_F) = \sum_{j=1}^d a_j H_F^{-1}(p_j).$$

We follow in part the treatment of simple L-functionals in Serfling (1980), §8.2.4.

Define

$$\Delta_{ln} = T_1(H_n) - T_1(H_F) - d T_1(H_F^{-1} H_n^{-1}).$$

Then

$$(3.4) \quad \Delta_{ln} = - \int_{-\infty}^{\infty} W_{H_n, H_F}(y) H_n(y) H_F(y) dy.$$

where we define

$$W_{G_1, G_2}(y) = \frac{K(G_1(y)) K(G_2(y))}{G_1(y) G_2(y)} - J(G_2(y)), \quad G_1(y), G_2(y),$$

L.s.m. G_1 is bounded and continuous a.e. Lebesgue and a.e. H_F^{-1} .

$$= 0, \quad G_1(y) G_2(y),$$

and

$$K(u) = \int_0^u J(t) dt.$$

From (3.4) we obtain two inequalities,

$$(3.5) \quad |\Delta_{ln}| \leq \|W_{H_n, H_F}\|_{L_1} \cdot \|H_n^{-1} H_F\|_{\infty},$$

and

$$(3.6) \quad |\Delta_{ln}| \leq \|W_{H_n, H_F}\|_{\infty} \cdot \|H_n^{-1} H_F\|_{L_1},$$

where $\|g\|_{\infty} = \sup_x |g(x)|$ and $\|g\|_{L_1} = \int |g(x)| dx$. We seek to establish

$$(3.7) \quad \sqrt{n} \Delta_{ln} \xrightarrow{P} 0$$

by analyzing the factors on the right-hand sides of (3.5) and (3.6).

LEMMA 3.1. Let J be as in Theorem 3.1. Then

$$\lim_{n \rightarrow \infty} \|G_1 - G_2\|_{\infty} = 0, \quad \|W_{G_1, G_2}\|_{L_1} = 0.$$

LEMMA 3.2. Let J be as in Theorem 3.2. Then

$$\lim_{n \rightarrow \infty} \|G_1 - G_2\|_{\infty} = 0, \quad \|W_{G_1, G_2}\|_{\infty} = 0.$$

(These are given as Lemmas 8.2.4A and 8.2.4E, respectively in Serfling (1980).)

LEMMA 3.3. If H_F is continuous, then

$$\|H_n - H_F\|_\infty = O_p(n^{-1}).$$

PROOF. Silverman (1976) establishes that the empirical stochastic process of a U-statistic array (indeed, of a more general type of array),

$$n^{\frac{1}{2}}(H_n(H_F^{-1}(t)) - t), \quad 0 \leq t \leq 1,$$

converges weakly in the Skorohod topology to an a.s. continuous Gaussian process, say W_t . By continuity of the mapping $\|\cdot\|_\infty$ with respect to the Skorohod topology, it follows that $n^{\frac{1}{2}}(H_n - H_F)\|_\infty \not\rightarrow \|W\|_\infty$ and hence that $n^{\frac{1}{2}}(H_n - H_F)\|_\infty = O_p(1)$. \square

LEMMA 3.4. Let H_F satisfy $|H_F'(1-H_F)| < \infty$. Let J be as in

Theorem 3.2. Then

$$E\|H_n - H_F\|_{L_1} = o(n^{-1}).$$

PROOF. Adapting the proof of Lemma 8.2.4D of Serfling (1980), write

$$H_n(y) - H_F(y) = (\frac{n}{m})^{-1} \int r_y(X_{i_1}, \dots, X_{i_m}),$$

with

$$r_y(\cdot) = I(h(\cdot) \leq y) - H_F(y).$$

Then

$$\|H_n - H_F\|_{L_1} = \|(\frac{n}{m})^{-1} \int r_y(X_{i_1}, \dots, X_{i_m})\| dy$$

and hence, using Tonelli's Theorem (Royden (1968), p. 270),

$$E\left(\|H_n - H_F\|_{L_1}\right) = E\left(\left(\frac{n}{m}\right)^{-1} \int r_y(X_{i_1}, \dots, X_{i_m})\right) dy.$$

Now, by a result on U-statistics (oeffding (1948); Serfling (1980), p. 183), and using Jensen's inequality, we have

$$E\left(\left(\frac{n}{m}\right)^{-1} \int r_y(X_{i_1}, \dots, X_{i_m})\right) \leq \left(\frac{n}{m}\right) E_y^2(X_1, \dots, X_m)^{\frac{1}{2}},$$

Thus

$$E\left(\|H_n - H_F\|_{L_1}\right) \leq n^{\frac{1}{2}} m^{-\frac{1}{2}} |H_F'(1-H_F)|. \quad \square$$

REMARK 3.1. In the proofs of Theorems 3.1 and 3.2, we will require (3.7).

Note that this follows from the conditions of Lemmas 3.1 and 3.3 together, as well as from the conditions of Lemmas 3.2 and 3.4 together. \square

Now, regarding Δ_{2n} , let us note that it may be written in the form

$$(3.8) \quad \Delta_{2n} = \sum_{j=1}^d a_j \hat{e}_{p_j, n}^{-1} p_j - \frac{p_i - H_n(p_i)}{H_F'(p_i)},$$

where \hat{e}_{p_j} denotes $H_F^{-1}(p_j)$ and $\hat{e}_{p_j, n}$ denotes $H_n^{-1}(p_j)$. In the case of simple L-statistics, the j-th term above is of the form

$$\hat{e}_{pn} - e_p - \frac{p - F(e_p)}{f(e_p)},$$

which is recognized to be the remainder term R_n in the Bahadur representation for sample quantiles (see Bahadur (1966) and Serfling (1980), p. 236.) Bahadur (1966) showed, under second-order differentiability conditions on F , that $R_n = o_p(n^{-3/4}(\log n)^{3/4})$. Ghosh (1971) showed $R_n = O_p(n^{-1/2})$ under only first-order differentiability conditions on F . The extension of Ghosh's result to the more general situation involving the terms in (3.8) is straightforward (details omitted), and we have

LEMMA 3.5. Let H_F have positive derivatives at its p_j -quantiles. Then $\eta_{2n} \geq 0$.

REMARKS 3.2. (1) The proofs of Theorems 3.1 and 3.2 are now straightforward, from Remark 3.1, Lemma 3.5, and the discussion at the beginning of this section.

(ii) To treat the alternative estimator $\tilde{T}(F_n)$, note that we need to deal with $\tilde{\Delta}_{1n} = \tilde{T}_1(F_n) - \tilde{T}_1(F) - d_1\tilde{T}_1(F, F_n, F)$, $i=1, 2$. Since $\tilde{T}_1(F) = T_1(H_F)$ and $\tilde{T}_1(F_n) = T_1(H_{F_n})$, we have the following analogues of (3.5) and (3.6):

$$(3.9) \quad |\tilde{\Delta}_{1n}| \leq \|u_{H_{F_n}, H_F}\|_{L_1} \cdot \|H_{F_n} - H_F\|_\infty$$

and

$$(3.10) \quad |\tilde{\Delta}_{2n}| \leq \|u_{H_{F_n}, H_F}\|_\infty \cdot \|H_{F_n} - H_F\|_{L_1}.$$

The proof (cf. analogues of Theorems 3.1 and 3.2) utilizes Lemmas 3.1, 3.2 and 3.5 without change, but requires analogues of Lemmas 3.3 and 3.4 with H_n replaced by H_{F_n} . Evidently these entail additional conditions on the kernel h . We shall not pursue these details here. \square

4. Complements.

4.1. Generalized M-statistics. An M-estimate (of location) may be defined in terms of the M-functional $T(\cdot)$ defined by

$$(4.1) \quad \int \psi(x - T(F)) dF(x) = 0,$$

where ψ is a given function. (See Huber (1977). For example.) Just as we

defined generalized L-functionals by replacing $T(F)$ by $T(H_F)$ for a specified L-functional $T(\cdot)$, we may define a generalized M-functional by putting H_F for F in (4.1). Thus a generalized M-statistic is given by $T(H_n)$. The analysis of such statistics follows standard lines with appropriate modifications due to the structure of H_n as a U-statistic.

4.2. Generalized R-statistics. Similar discussion.

4.3. Berry-Esséen theorems for generalized L-statistics. A method used in Serfling (1980), pp. 287-290, for simple L-statistics can in principle be extended to generalized L-statistics. First, a higher-order differential for GL-functionals is needed; this can be obtained by continuing to differentiate in (2.12) or (2.13). Secondly, certain moment theorems for $\|H_n - H_F\|_{L_2}$ are needed; these are straightforward.

4.4. Strong convergence of generalized L-statistics. For simple L-statistics, a law of iterated logarithm analogue of Theorem 3.1 follows under second-order differentiability of F at its p_j -quantiles. The argument (see Serfling (1980), p. 281) uses LL results for $\|F_n - F\|_\infty$ and for the remainder term in Bahadur's representation for sample quantiles. With certain generalizations of these results, the argument should extend to GL-statistics.

REFERENCES

- Bahadur, R. R. (1966). "A note on quantiles in large samples," Ann. Math. Statist., 37, 577-580.

Unclassified
SECURITY CLASSIFICATION OF THIS PAGE

REPORT DOCUMENTATION PAGE	
1. REPORT NUMBER	2. GOVT ACCESSION NO.
ONR No. 81-8	3. RECIPIENT CATALOG NUMBER
4. TITLE GENERALIZED L-, M- AND R-STATISTICS	5. TYPE OF REPORT & PERIOD COVERED Technical Report
5. CONTRACT OR GRANT NUMBER(s) ONR No. N00014-79-C-0801	6. PERFORMING ORGANIZATION REPORT NO. Technical Report No. 149
7. AUTHOR(s) Robert J. Serfling	8. CONTRACT OR GRANT NUMBER(s)
9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Mathematical Sciences The Johns Hopkins University Baltimore, Maryland 21218	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
11. CONTROLLING OFFICE NAME & ADDRESS Office of Naval Research Statistics and Probability Program Arlington, Virginia 22217	12. REPORT DATE August 1, 1981
13. NUMBER OF PAGES 20	14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office) Unclassified
15. SECURITY CLASS (of this report) Unclassified	16. DISTRIBUTION STATEMENT (of this report) Approved for public release; distribution unlimited.
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from report) Approved for public release; distribution unlimited.	18. SUPPLEMENTARY NOTES in <u>New Directions and Recent Results in Algorithms and Complexity</u> , ed. by J. F. Traub, Academic Press.
19. KEY WORDS Hodges-Lehmann estimator; asymptotic normality.	20. ABSTRACT A class of statistics generalizing U-statistics and L-statistics, and containing other varieties of statistic as well, such as trimmed U-statistics, is studied. Using the differentiable statistical function approach, differential approximations are obtained and the influence curves of these generalized L-statistics are derived. These results are employed to establish asymptotic normality for such statistics. Parallel generalizations of M- and R-statistics are noted. Strong convergence, Berry-Esséen rates, and computational aspects are discussed.
von Mises, R. (1947). "Limit theorems for dissociated random variables," <u>Adv. Appl. Prob.</u> , <u>8</u> , 806-819.	
Bickel, P. J. and Lehmann, E. L. (1976). "Descriptive statistics for nonparametric models. III. Dispersion," <u>Ann. Statist.</u> , <u>4</u> , 1139-1158.	
Blum, M., Floyd, R. W., Pratt, V., Rivest, R., Tarjan, R. E. (1973). "Time bounds for selection," <u>J. Comp. and System Sci.</u> , <u>7</u> , 448-461.	
Floyd, R. W. and Rivest, R. (1975). "Expected time bounds for selection," <u>Commun. A.C.M.</u> , <u>18</u> , 165-172.	
Ghosh, J. K. (1971). "A new proof of the Bahadur representation of quantiles and an application," <u>Ann. Math. Statist.</u> , <u>42</u> , 1957-1961.	
Hampel, F. R. (1968). "Contributions to the theory of robust estimation," Ph.D. dissertation, Univ. of California-Berkeley.	
Hampel, F. R. (1974). "The influence curve and its role in robust estimation," <u>J. Amer. Statist. Assoc.</u> , <u>69</u> , 383-397.	
Hoeffding, W. (1948). "A class of statistics with asymptotically normal distribution," <u>Ann. Math. Statist.</u> , <u>19</u> , 293-325.	
Huber, P. J. (1977). <u>Robust Statistical Procedures</u> , SIAM, Philadelphia.	
Royden, H. L. (1968). <u>Real Analysis</u> , 2nd. ed., Macmillan, New York.	
Serfling, R. J. (1980). <u>Approximation Theorems of Mathematical Statistics</u> , Wiley, New York.	
Shapiro, M. J. (1976). "Geometry and Statistics: problems at the interface," in <u>New Directions and Recent Results in Algorithms and Complexity</u> , ed. by J. F. Traub, Academic Press.	

END
DATE
FILMED

10-81.

DTIC